

# Online Neural Network-based Language Identification

Master's Thesis Daniel H. Draper

INTERACTIVE SYSTEMS LAB - INSTITUTE FOR ANTHROPOMATICS AND ROBOTICS



## **Contents**



- Motivation
- Related Work
- Theory of LID
- **Experimental Setup**
- Demo
- Post Processing
- Results



## **Motivation**



- ASR and SLT commercially great success ["OK, Google", Siri, Jibbigo]
- Best performing ASR models are still trained on one source language.[1]
- Performance increase for multilingual tasks[1], UI streamlining.



**Experimental Setup** 

## Motivation



- ASR and SLT commercially great success ["OK, Google", Siri, Jibbigo]
- Best performing ASR models are still trained on one source language.[1]
- Performance increase for multilingual tasks[1], UI



**Experimental Setup** 

## Motivation



- ASR and SLT commercially great success ["OK, Google", Siri, Jibbigo]
- Best performing ASR models are still trained on one source language.[1]
- Performance increase for multilingual tasks[1], UI streamlining.





## Lecture Translator

- Low latency, online application



- Many languages challenging, great potential.





## Lecture Translator

- Low latency, online application
- Small source language amount make performance increase likely



- SLT already employed consistently[2]
- Many languages challenging, great potential.





## Lecture Translator

- Low latency, online application
- Small source language amount make performance increase likely

## European Parliament

- SLT already employed consistently[2]
- Many languages challenging, great potential





## Lecture Translator

- Low latency, online application
- Small source language amount make performance increase likely

## European Parliament

- SLT already employed consistently[2]
- Many languages challenging, great potential.



## **Related Work I**



#### Historically

- PRLM<sup>a</sup>, PPRM<sup>b</sup>[3]
- Gaussian Mixture Models (GMM)s [4]

<sup>a</sup>single-language phone recognition followed by language-dependent, interpolated n-gram language modeling

<sup>b</sup>multiple single-language phone recognizers and language-dependent parallel phone recognition



## **Related Work I**



#### Historically

- PRLM<sup>a</sup>, PPRM<sup>b</sup>[3]
- Gaussian Mixture Models (GMM)s [4]

<sup>a</sup>single-language phone recognition followed by language-dependent, interpolated n-gram language modeling

<sup>b</sup>multiple single-language phone recognizers and language-dependent parallel phone recognition



## Related Work II



## Modern Approaches

- VSM<sup>a</sup> + PPRM[5]
- <sup>a</sup>Vector Space Modelling



## Related Work II



## Modern Approaches

- VSM<sup>a</sup> + PPRM[5]
- I-Vector approaches[6][7]

<sup>a</sup>Vector Space Modelling



Demo

## Related Work III/Based on



## Similar Approaches

- $\blacksquare$  Matejka et. al[8] similar setup: BNF  $\to$  LID, with averaging 5 LID nets.
- [9] M. Heck et al. evaluate LID approaches prospect of Lecture Translator integration: PPRM/PRLM + Hybrid.
- Experimental baseline setup from Markus [10]: "Language Adaptive DNNs for Improved Low Resource Speech Recognition"

**Experimental Setup** 

Theory of LID

## Related Work III/Based on



## Similar Approaches

- Matejka et. al[8] similar setup: BNF → LID, with averaging 5 LID nets.
- [9] M. Heck et al. evaluate LID approaches prospect of Lecture Translator integration: PPRM/PRLM + Hybrid.
- Experimental baseline setup from Markus [10]: "Language



## Related Work III/Based on



#### Similar Approaches

- Matejka et. al[8] similar setup: BNF → LID, with averaging 5 LID nets.
- [9] M. Heck et al. evaluate LID approaches prospect of Lecture Translator integration: PPRM/PRLM + Hybrid.

 Experimental baseline setup from Markus [10]: "Language Adaptive DNNs for Improved Low Resource Speech Recognition"





• Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$ 

$$L^O = \arg\max_l p(S|L_l)$$

$$L^O = \operatorname{arg\,max}_I P(\Upsilon|L_I)$$

With Viterbi decoding on set of phone models *M*:

$$\Upsilon^O = \operatorname{arg\,max}_v P(S|v, M)$$

$$L^O = \operatorname{arg\,max}_I \sum_{\forall \Upsilon} P(S|\Upsilon, M) P(\Upsilon|L_I)$$





• Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$ 

## Formulation of LID Task[11]

$$L^O = \operatorname{arg\,max}_I p(S|L_I)$$

With segmentation into phones v (sequence  $\Upsilon$ ):

$$L^O = \arg\max_l P(\Upsilon|L_l)$$

$$\Upsilon^{O} = \operatorname{arg\,max}_{v} P(S|v, M)$$

$$L^O = \operatorname{arg\,max}_I \sum_{\forall \Upsilon} P(S|\Upsilon, M) P(\Upsilon|L_I)$$



Demo

Results



• Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$ 

## Formulation of LID Task[11]

$$L^O = \operatorname{arg\,max}_I p(S|L_I)$$

With segmentation into phones v (sequence  $\Upsilon$ ):

$$L^O = \operatorname{arg\,max}_I P(\Upsilon|L_I)$$

With Viterbi decoding on set of phone models *M*:

$$\Upsilon^O = \operatorname{arg\,max}_{\upsilon} P(S|\upsilon, M)$$

$$L^{O} = \operatorname{arg\,max}_{I} \sum_{\forall \Upsilon} P(S|\Upsilon, M) P(\Upsilon|L_{I})$$





• Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$ 

## Formulation of LID Task[11]

$$L^O = \arg\max_I p(S|L_I)$$

With segmentation into phones v (sequence  $\Upsilon$ ):

$$L^O = \operatorname{arg\,max}_I P(\Upsilon|L_I)$$

With Viterbi decoding on set of phone models *M*:

$$\Upsilon^{\mathcal{O}} = \operatorname{arg\,max}_{v} P(\mathcal{S}|v, M)$$

Combining last 2 equations:

$$L^{O} = \operatorname{arg\,max}_{I} \sum_{\forall \Upsilon} P(S|\Upsilon, M) P(\Upsilon|L_{I})$$





• Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$ 

## Formulation of LID Task[11]

$$L^O = \operatorname{arg\,max}_I p(S|L_I)$$

With segmentation into phones v (sequence  $\Upsilon$ ):

$$L^O = \operatorname{arg\,max}_I P(\Upsilon|L_I)$$

With Viterbi decoding on set of phone models *M*:

$$\Upsilon^{O} = \operatorname{arg\,max}_{v} P(S|v, M)$$

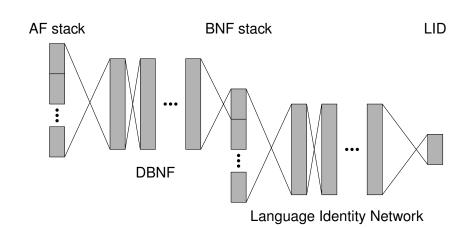
Combining last 2 equations:

$$L^{O} = \operatorname{arg\,max}_{I} \sum_{orall \Upsilon} P(S|\Upsilon, M) P(\Upsilon|L_{I})$$



## **Experimental Setup**







#### **Data Sets**



#### Euronews 2014

- 10 Languages (EN, FR, DE, AR, ES, IT, PO, PT, RU, TR)
- Original Corpus: 72h / Language, Reduced corpus: 18h / Language (Random Sampling of 10.000 speakers)
- 80 % train data, 10 % each dev/test set.

#### Lecture Data

- 3 Languages (EN, FR, DE) 10h per language.
- KIT lectures, InterACT25, DGA talks
- Evaluation of EUNews results, Proof-of-concept for Lecture Translator integration.



#### **Data Sets**



#### Euronews 2014

- 10 Languages (EN, FR, DE, AR, ES, IT, PO, PT, RU, TR)
- Original Corpus: 72h / Language, Reduced corpus: 18h / Language (Random Sampling of 10.000 speakers)
- 80 % train data, 10 % each dev/test set.

#### Lecture Data

- 3 Languages (EN, FR, DE) 10h per language.
- KIT lectures, InterACT25, DGA talks
- Evaluation of EUNews results, Proof-of-concept for Lecture Translator integration.



#### Data Sets II



#### European Parliament

- 7 Languages (EN, FR, DE, ES, IT, PO, PT) 3.6h per language
- Simultaneous translation of all EUParliament speeches
- Further Evaluation, Proof-of-concept for EUParl integration.



## **Feature Extraction**



990 Demo Post Processing Results

Experimental Setup

000000

Related Work Theory of LID Motivation Daniel H. Draper - Online Neural Network-based Language Identification

## **DBNF**



990 Post Processing Results Demo

Related Work Daniel H. Draper - Online Neural Network-based Language Identification

Theory of LID

Experimental Setup

000000

Motivation

## **DNN Structure**



990 Results Demo Post Processing

Experimental Setup

000000

Motivation Related Work Theory of LID Daniel H. Draper - Online Neural Network-based Language Identification

## **Demo**



《ロトイラトモミトモミト 夏 シュで Motivation Related Work Theory of LID Experimental Setup **Demo** Post Processing Results

Daniel H. Draper – Online Neural Network-based Language Identification

## **Evaluation Metrics**



 Image: Control of the control of

16/24

Motivation Related Work Theory of LID Experimental Setup Demo Post Processing

O Daniel H. Draper - Online Neural Network-based Language Identification

## **Smoothing Filters**





Motivation Related Work Theory of LID Experimental Setup

Daniel H. Draper - Online Neural Network-based Language Identification

## Results



990 Post Processing Results Demo

Experimental Setup

Motivation Related Work Theory of LID Daniel H. Draper - Online Neural Network-based Language Identification

## **Future Work**



## Questions?

990

Daniel H. Draper - Online Neural Network-based Language Identification

Theory of LID

Related Work

Motivation

Experimental Setup

Demo

Post Processing

20/24

Results

## References I



- Z. Tang, L. Li, and D. Wang, "Multi-task recurrent model for true multilingual speech recognition," arXiv preprint arXiv:1609.08337, 2016.
- "Technology and corpora for speech to speech translation," http://tcstar.org/.
- M. A. Zissman et al., "Comparison of four approaches to automatic language identification of telephone speech," IEEE Transactions on speech and audio processing, vol. 4, no. 1, p. 31, 1996.
- P. A. Torres-Carrasquillo, E. Singer, M. A. Kohler, R. J. Greene, D. A. Reynolds, and J. R. Deller Jr, "Approaches to language identification using gaussian mixture models and shifted delta cepstral features." in Interspeech, 2002.



## References II



- H. Li, B. Ma, and C. H. Lee, "A vector space modeling approach to spoken language identification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 1, pp. 271–284, Jan 2007.
- Y. Song, B. Jiang, Y. Bao, S. Wei, and L. R. Dai, "I-vector representation based on bottleneck features for language identification," *Electronics Letters*, vol. 49, no. 24, pp. 1569–1570, November 2013.
- L. F. D'haro Enríquez, O. Glembek, O. Plchot, P. Matějka, M. Soufifar, R. d. Córdoba Herralde, and J. Černockỳ, "Phonotactic language recognition using i-vectors and phoneme posteriogram counts," 2012.



## References III



- P. Matejka, L. Zhang, T. Ng, H. S. Mallidi, O. Glembek, J. Ma, and B. Zhang, "Neural network bottleneck features for language identification," Proc. IEEE Odyssev, pp. 299-304, 2014.
- M. Heck, S. StÃ<sup>1</sup>/<sub>4</sub>ker, and A. Waibel, "A hybrid phonotactic language identification system with an svm back-end for simultaneous lecture translation." in 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), March 2012, pp. 4857–4860.



#### References IV



- M. Mà <sup>1</sup>/<sub>4</sub>ller, S. Stà <sup>1</sup>/<sub>4</sub>ker, and A. Waibel, "Language adaptive dnns for improved low resource speech recognition," in Proceedings of the 17th Annual Conference of the International Speech Communication Association (INTERSPEECH), San Francisco, USA, September 8-12 2016.
- H. Li, B. Ma, and K. A. Lee, "Spoken language recognition: From fundamentals to practice," *Proceedings of the IEEE*, vol. 101, no. 5, pp. 1136–1159, May 2013.